

Text and spatial data mining

Finn Årup Nielsen

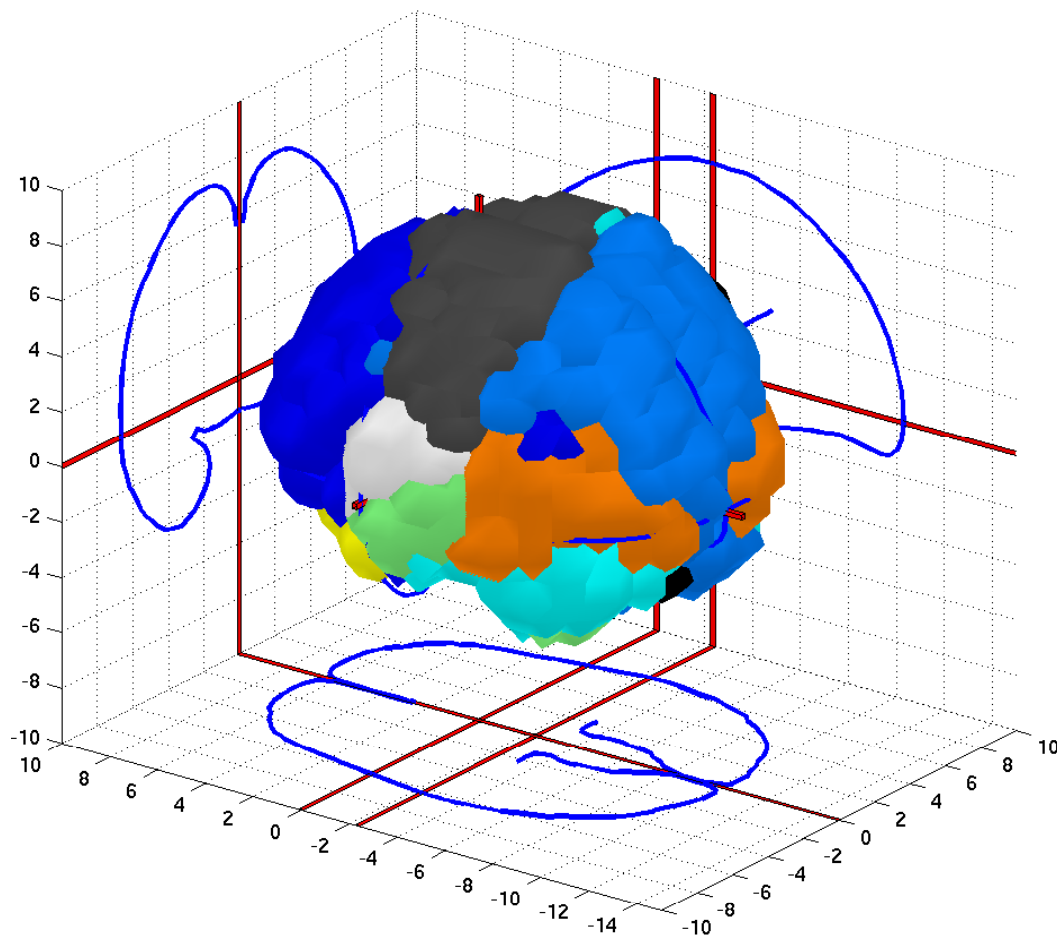
Lundbeck Foundation Center for Integrated Molecular Brain Imaging
at

Neurobiology Research Unit
Copenhagen University Hospital Rigshospitalet
and

Informatics and Mathematical Modelling
Technical University of Denmark

September 29, 2006

Parcellation of the human brain



Parcellation of the human brain by combining text mining and spatial data mining within a neuroinformatics database.

Text mining: Analysis of scientific abstracts.

Spatial data mining: Modeling of the distribution of Talairach coordinates.

Seek communality between the the text representation and spatial representation by multivariate analysis.

Brede Database

WOBIB: 27 - Epstein, Kanwisher (1998) A cortical repres ...

Bib -> [Asymmetry](#) | [Author](#) | [ICA](#) | [NMF](#) | [Novelty](#) | [Statistics](#) | [SVD](#) | [Title](#) | [WOBIB](#) |

Exp -> [Alphabetic](#) | [Asymmetry](#) | [ICA](#) | [NMF](#) | [Novelty](#) | [SVD](#) | [WOEXP](#) | [WOEXT](#) |

Ext -> [Alphabetic index](#) | [Map](#) | [Roots](#) | [\[Brede \]](#) | Loc -> [Statistics](#) |

R. Epstein; N. Kanwisher. [A cortical representation of the local visual environment](#). *Nature* **392**(6676):598-601, 1998. PMID: [9560155](#). DOI: [10.1038/33402](#). WOBIB: [27](#).

Medial temporal brain regions such as the hippocampal formation and parahippocampal cortex have been generally implicated in navigation and visual memory. However, the specific function of each of these regions is not yet clear. Here we present evidence that a particular area within human parahippocampal cortex is involved in a critical component of navigation: perceiving the local visual environment. This region, which we name the 'parahippocampal place area' (PPA), responds selectively and automatically in functional magnetic resonance imaging (fMRI) to passively viewed scenes, but only weakly to single objects and not at all to faces. The critical factor for this activation appears to be the presence in the stimulus of information about the layout of local space. The response in the PPA to scenes with spatial layout but no discrete objects (empty rooms) is as strong as the response to complex meaningful scenes containing multiple objects (the same rooms furnished) and over twice as strong as the response to arrays of multiple objects without three-dimensional spatial context (the furniture from these rooms on a blank background). This response is reduced if the surfaces in the scene are rearranged so that they no longer define a coherent space. We propose that the PPA represents places by encoding the geometry of the local environment.

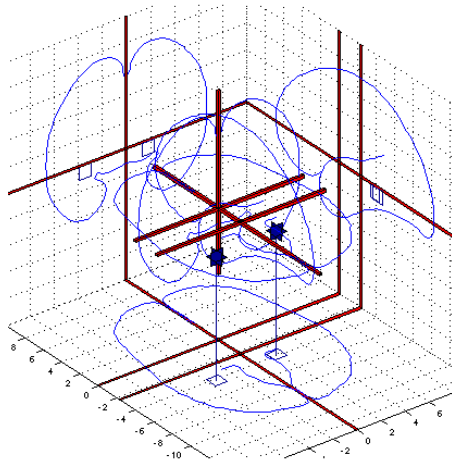


Figure 1: Web-page generated for the paper (Epstein and Kanwisher, 1998) with abstract and Talairach coordinates displayed in a corner cube environment (Rehm et al., 1998).

Brede Database (Nielsen, 2003) inspired by the BrainMap database (Fox and Lancaster, 1994).

Presently 186 papers and the information used is:

Abstract of the article

3D coordinates representing change in brain activation, grey matter variation or site of lesion: so-called Talairach coordinates (Talairach and Tournoux, 1988).

Non-negative matrix factorization

Multivariate analysis: Non-negative matrix factorization (NMF) decomposes a non-negative data matrix $\mathbf{X}(N \times P)$ (Lee and Seung, 1999)

$$\mathbf{X} = \mathbf{WH} + \mathbf{U}, \quad (1)$$

where $\mathbf{W}(N \times K)$ and $\mathbf{H}(K \times P)$ are also non-negative matrices.

“Euclidean” cost function for

$$E_{\text{“eucl”}} = \|\mathbf{X} - \mathbf{WH}\|_F^2 \quad (2)$$

Iterative algorithm (Lee and Seung, 2001)

$$\mathbf{H}_{kp} \leftarrow \mathbf{H}_{kp} \frac{(\mathbf{W}^\top \mathbf{X})_{kp}}{(\mathbf{W}^\top \mathbf{WH})_{kp}} \quad (3)$$

$$\mathbf{W}_{nk} \leftarrow \mathbf{W}_{nk} \frac{(\mathbf{XH}^\top)_{nk}}{(\mathbf{WHH}^\top)_{nk}}. \quad (4)$$

Partial least squares

One of the variations of partial least squares: Singular value decomposition of a inner product matrix (McIntosh et al., 1996)

$$\mathbf{ULV}^T = \text{svd}(\mathbf{X}^T \mathbf{Y}) \quad (5)$$

Probably most suitable for data that is symmetric, i.e., both positive and negative.

“Non-negative partial least squares”

$$\mathbf{WH} = \text{nmf}(\mathbf{X}^T \mathbf{Y}) \quad (6)$$

Should get two non-negative matrices (\mathbf{X} and \mathbf{Y}) as input.

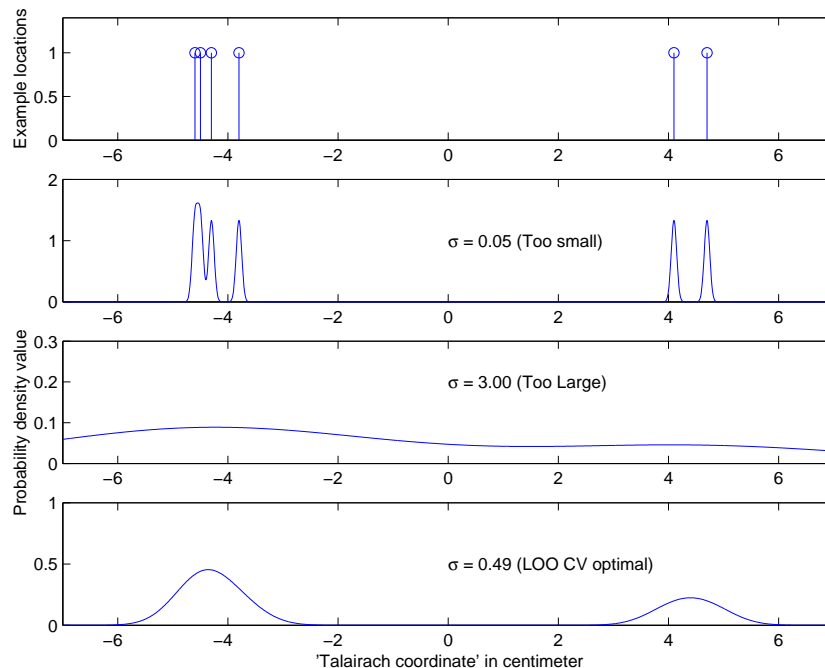
First matrix: Bag-of-words matrix

	'memory'	'visual'	'motor'	'time'	'retrieval'	...
Fujii	6	0	1	0	4	...
Maddock	5	0	0	0	0	...
Tsukiura	0	0	4	0	0	...
Belin	0	0	0	0	0	...
Ellerman	0	0	0	5	0	...
⋮	⋮	⋮	⋮	⋮	⋮	⋮

Representation of the abstracts of the papers in a bag-of-words matrix: (abstract \times words)-matrix: Each element counts of the frequency of a word occurring in an abstract text (Salton et al., 1975).

Exclusion of stop words: common words, brain anatomy, ... Mostly words for brain function left (Nielsen et al., 2005).

Second matrix: Voxelization matrix



Regard the Talairach coordinates in an article as being generated from a distribution $p(\mathbf{z})$, where \mathbf{z} is in 3D Talairach space (Fox et al., 1997).

Kernel methods (L kernels centered on each location: μ_l) with homogeneous Gaussian kernel in 3D Talairach space \mathbf{z}

$$\hat{p}(\mathbf{z}) = \frac{(2\pi\sigma^2)^{-3/2}}{L} \sum_l e^{-\frac{1}{2\sigma^2}(\mathbf{z}-\mu_l)^2}$$

σ^2 fixed ($\sigma = 1\text{cm}$) or optimized with leave-one-out cross-validation (Nielsen and Hansen, 2002).

Details

Coarse sampling of the volume with 8mm voxels: $\hat{p}(\mathbf{z}) \equiv \mathbf{y}$. This \mathbf{y}

Corresponds to the Talairach coordinates in one scientific article and when stacked: $\mathbf{Y}(N \times P)$

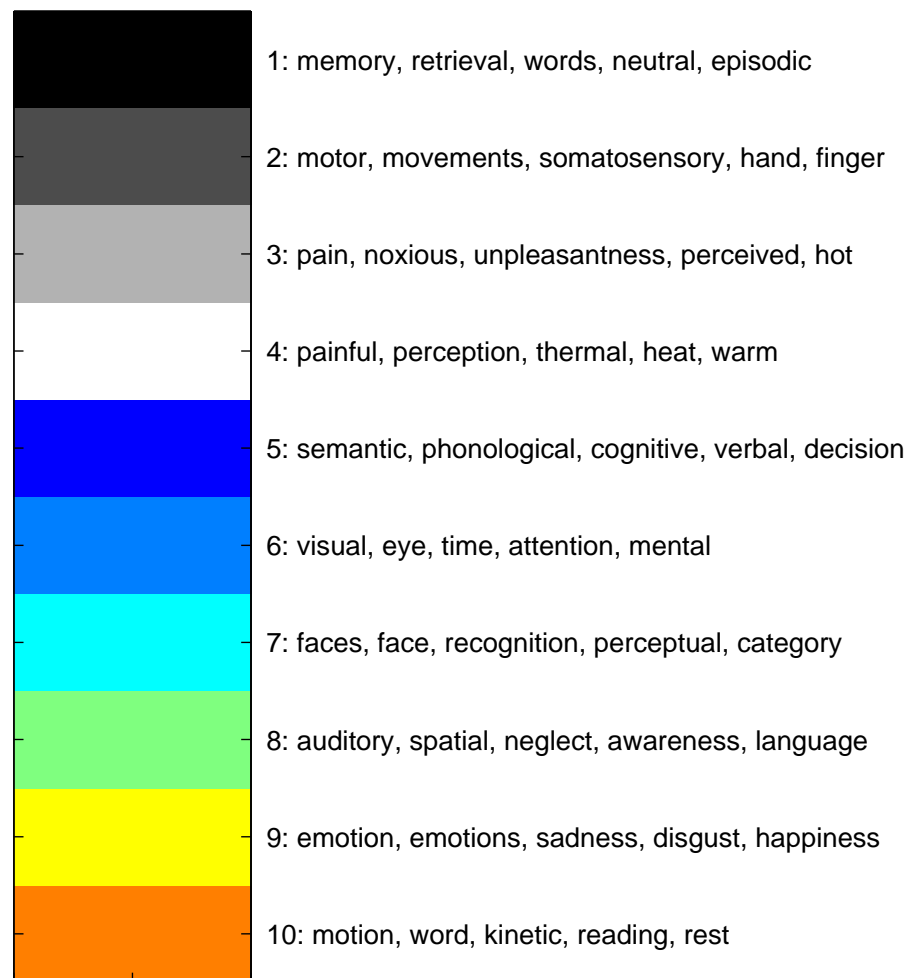
Restriction to gray matter regions using a anatomically labeled brain “AAL” (Tzourio-Mazoyer et al., 2002).

Number of components in NMF: Using a rule of thumb we select $K = \sqrt{N/2}$ (Mardia et al., 1979): $K = 10 \approx \sqrt{186/2}$.

Several runs of NMF with 50'000 iterations each, and with different initialization each time.

Exclusive assignment: Winner-takes-all function on \mathbf{W} and \mathbf{H} .

Resulting NMF components



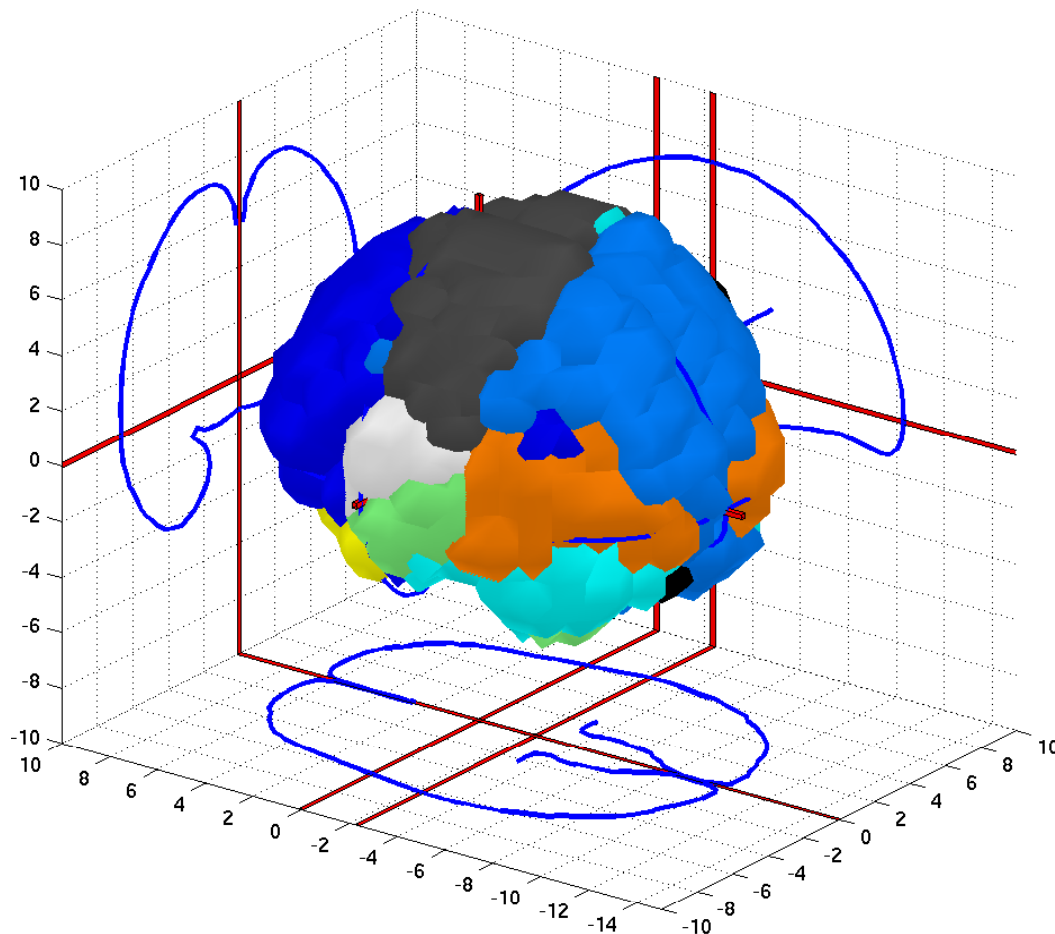
Loadings on words in W

$W(P \times K) = W(470 \times 10)$, i.e., 10 components (“functions”) and 470 words.

Most dominant functions listed at the top: memory, motor, pain, . . .

5 most loaded words listed for each function.

Dorsolateral surface view



\mathbf{H} in Talairach space: $\mathbf{H}(K \times Q) = \mathbf{H}(10 \times 2492)$, i.e., 10 components and 2492 voxels.

Occipital and parietal lobe: “visual”, “eye”.

Central sulcus: “Motor”, “movements”, “somatosensory”.

Temporal cortex: “auditory”, “spatial”, “neglect”

Medial view

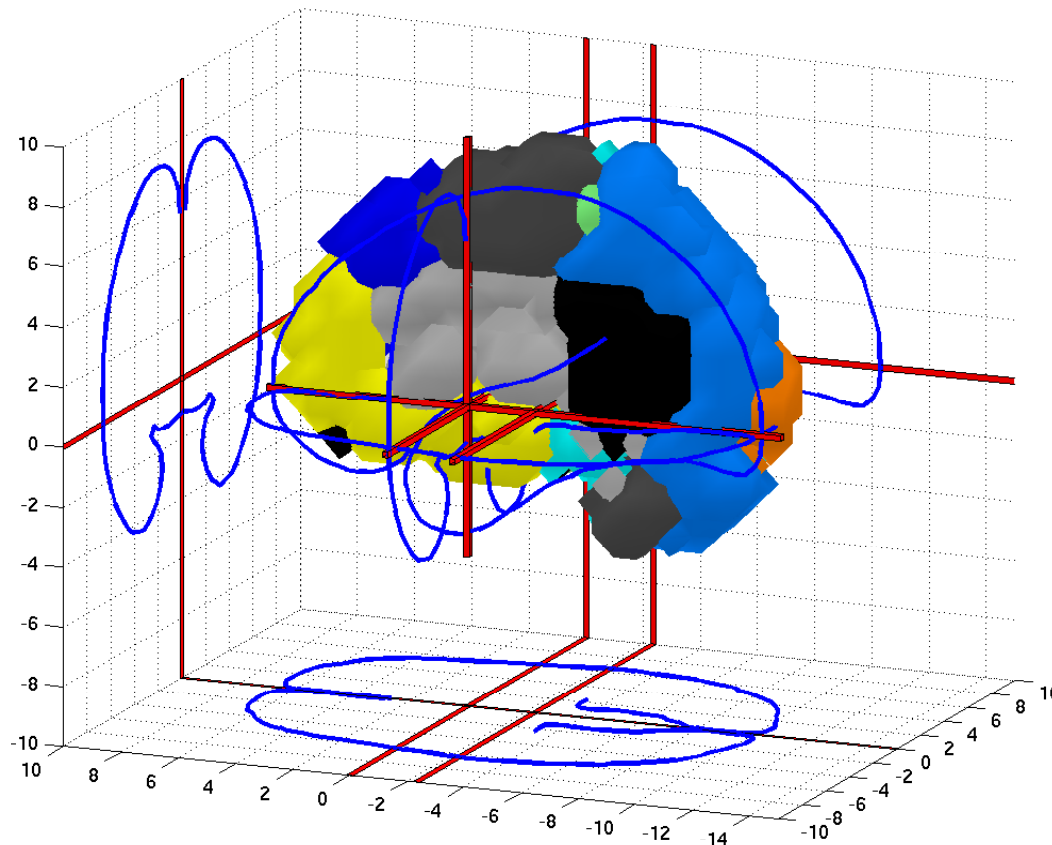


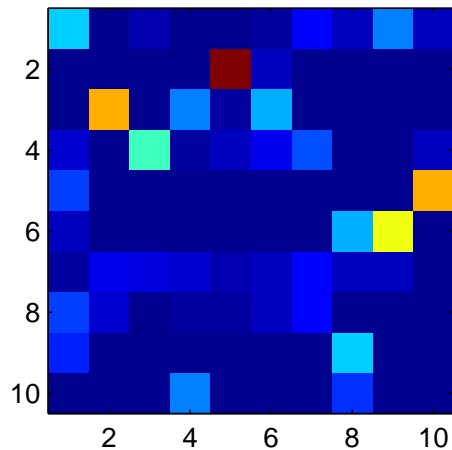
Figure 2: Medial surface view of the labeled right hemisphere. Seen from the left

“Memory” in posterior cingulate area. Probably due to the many articles about memory and the posterior cingulate in the Brede database. Episodic memory retrieval is associated with posterior cingulate (Cabeza and Nyberg, 2000).

“Emotion” in the medial frontal area, e.g., amygdala.

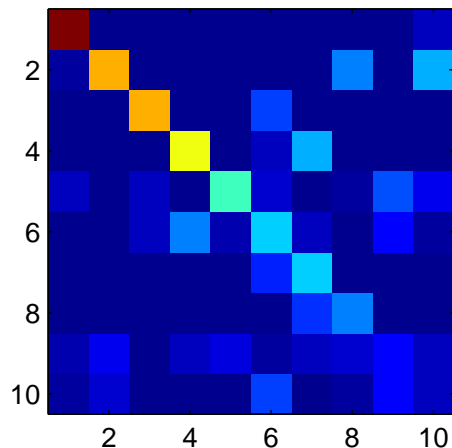
“Pain” in anterior cingulate, thalamus, insula. Previously noted: (Ingvar, 1999).

How stable are the results?



The parcellation varies between runs of the NMF and when different parts of the data set are used.

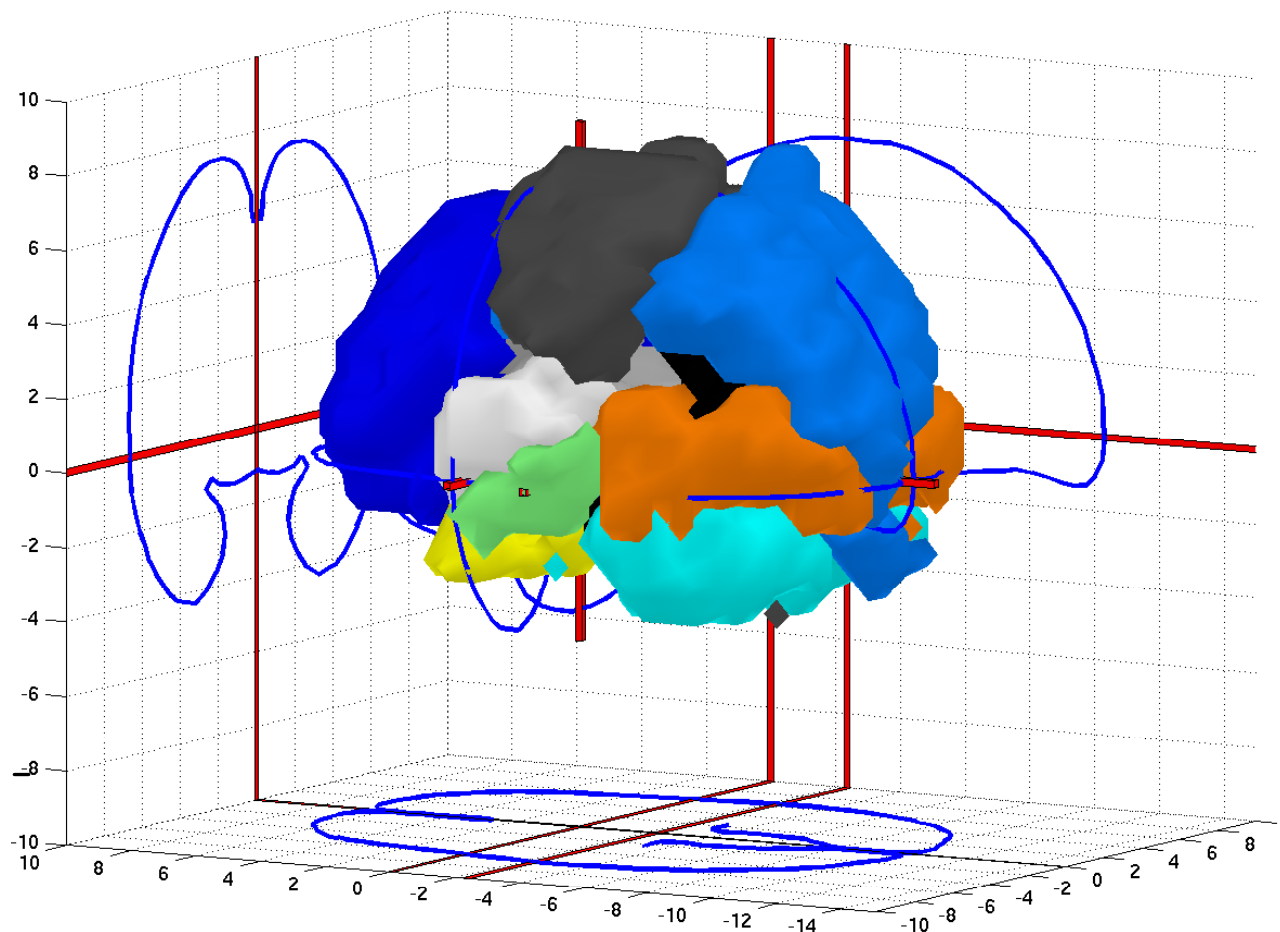
The different components are not matched between runs, e.g., the 2nd component in the first run might match the 5rd component in the second run.



It is possible to match the components, e.g., with the algorithm suggested by (Meilā, 2002) or by the “Hungarian method” (Roth et al., 2002).

Example “confusion matrices” (\mathbf{HH}^T) in the figure appear with one instance of half-split resampling between the all papers before and after sorting.

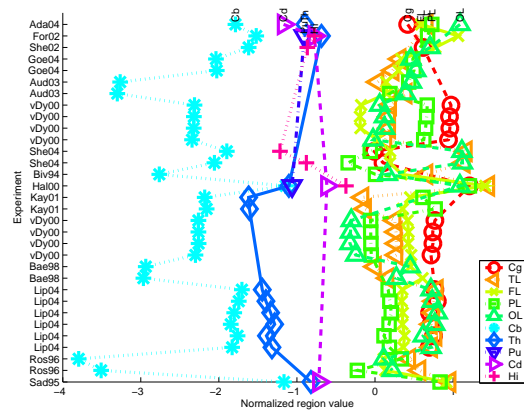
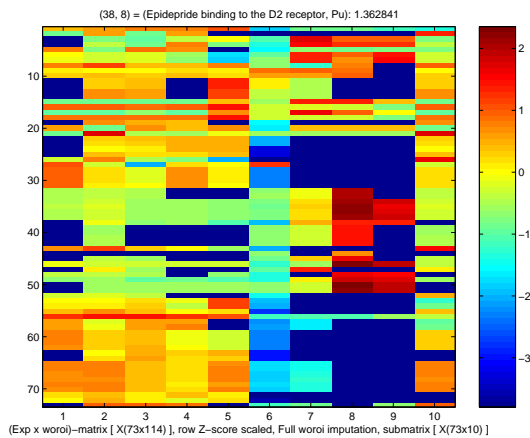
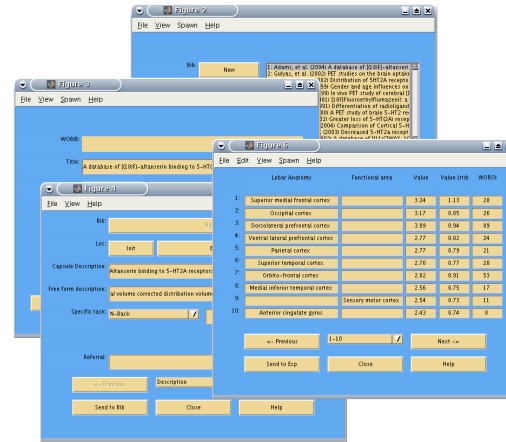
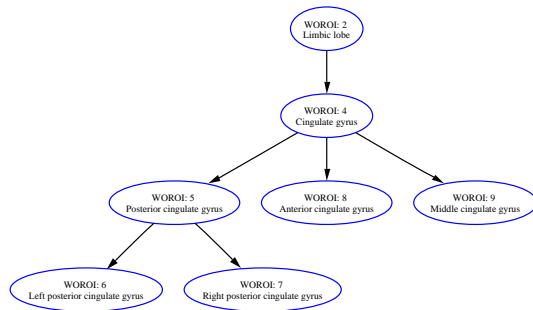
Consistent parcellation



Parcellation after masking those voxels that are consistently parcellated.

Threshold determined by a permutation test and matching by the Hungarian method.

... and molecular neuroimaging?



Few molecular neuroimaging studies with Talairach coordinates — most are based on analysis in pre-defined regions.

Building a taxonomy for brain regions.

Match regions in papers
with taxonomy entries

Build algorithms that are able to handle missing data and different scalings.

Summary

It is possible to automatically perform a high-level coarse parcellation of the entire human brain.

The results appear in accordance with general consensus in human brain mapping.

We use the Brede Database and rely only on abstract and Talairach coordinate information.

We suggest “non-negative partial least squares” as a combination of non-negative matrix factorization and partial least squares.

References

- Cabeza, R. and Nyberg, L. (2000). Imaging cognition II: An empirical review of 275 PET and fMRI studies. *Journal of Cognitive Neuroscience*, 12(1):1–47. PMID: 10769304. <http://jocn.mitpress.org/cgi/content/abstract/12/1/1>.
- Epstein, R. and Kanwisher, N. (1998). A cortical representation of the local visual environment. *Nature*, 392(6676):598–601. PMID: 9560155. DOI: 10.1038/33402. ISSN 0028-0836.
- Fox, P. T. and Lancaster, J. L. (1994). Neuroscience on the net. *Science*, 266(5187):994–996. PMID: 7973682.
- Fox, P. T., Lancaster, J. L., Parsons, L. M., Xiong, J.-H., and Zamarripa, F. (1997). Functional volumes modeling: Theory and preliminary assessment. *Human Brain Mapping*, 5(4):306–311. <http://www3.interscience.wiley.com/cgi-bin/abstract/56435/START>.
- Ingvar, M. (1999). Pain and functional imaging. *Philosophical Transactions of the Royal Society of London. Series B, Biological Sciences*, 354(1387):1347–1358. PMID: 10466155.
- Lee, D. D. and Seung, H. S. (1999). Learning the parts of objects by non-negative matrix factorization. *Nature*, 401(6755):788–791. PMID: 10548103.
- Lee, D. D. and Seung, H. S. (2001). Algorithms for non-negative matrix factorization. In Leen, T. K., Dietterich, T. G., and Tresp, V., editors, *Advances in Neural Information Processing Systems 13: Proceedings of the 2000 Conference*, pages 556–562, Cambridge, Massachusetts. MIT Press. <http://hebb.mit.edu/people/seung/papers/nmfconverge.pdf>. CiteSeer: <http://citeseer.ist.psu.edu/-lee00algorithms.html>.
- Mardia, K. V., Kent, J. T., and Bibby, J. M. (1979). *Multivariate Analysis*. Probability and Mathematical Statistics. Academic Press, London. ISBN 0124712525.
- McIntosh, A. R., Bookstein, F. L., Haxby, J. V., and Grady, C. L. (1996). Spatial pattern analysis of functional brain images using Partial Least Square. *NeuroImage*, 3(3 part 1):143–157. PMID: 9345485. ftp://ftp.rotman-baycrest.on.ca/pub/Randy/PLS/pls_article.pdf.

- Meilā, M. (2002). Comparing clusterings. Technical Report 418, Department of Statistics, University of Washington, Seattle, Washington. <http://www.stat.washington.edu/www/research/reports/2002-tr418.ps>.
- Nielsen, F. Å. (2003). The Brede database: a small database for functional neuroimaging. *NeuroImage*, 19(2). <http://208.164.121.55/hbm2003/abstract/abstract906.htm>. Presented at the 9th International Conference on Functional Mapping of the Human Brain, June 19–22, 2003, New York, NY. Available on CD-Rom.
- Nielsen, F. Å., Balslev, D., and Hansen, L. K. (2005). Mining the posterior cingulate: Segregation between memory and pain component. *NeuroImage*, 27(3):520–532. DOI: 10.1016/j.neuroimage.2005.04.034.
- Nielsen, F. Å. and Hansen, L. K. (2002). Modeling of activation data in the BrainMap™ database: Detection of outliers. *Human Brain Mapping*, 15(3):146–156. DOI: 10.1002/hbm.10012. <http://www3.interscience.wiley.com/cgi-bin/abstract/89013001/>. CiteSeer: <http://citeseer.ist.psu.edu/nielsen02modeling.html>.
- Rehm, K., Lakshminarayan, K., Frutiger, S. A., Schaper, K. A., Sumners, D. L., Strother, S. C., Anderson, J. R., and Rottenberg, D. A. (1998). A symbolic environment for visualizing activated foci in functional neuroimaging datasets. *Medical Image Analysis*, 2(3):215–226. PMID: 9873900. <http://www.sciencedirect.com/science/article/B6W6Y-45PJY0D-7/1-/48196224354fdd62ea8c5a0d85379b07>.
- Roth, V., Lange, T., Braun, M., and Buhmann, J. (2002). A resampling approach to cluster validation. In Härdle, W. and Rönz, B., editors, *Proceedings in Computational Statistics: 15th Symposium Held in Berlin, Germany 2002 (COMPSTAT2002)*, pages 123–128, Heidelberg, Germany. Physica-Verlag. <http://www.ml.inf.ethz.ch/publications/roth.compstat02.pdf>.
- Salton, G., Wong, A., and Yang, C. S. (1975). A vector space model for automatic indexing. *Communication of the ACM*, 18:613–620.
- Talairach, J. and Tournoux, P. (1988). *Co-planar Stereotaxic Atlas of the Human Brain*. Thieme Medical Publisher Inc, New York. ISBN 0865772932.
- Tzourio-Mazoyer, N., Landeau, B., Papathanassiou, D., Crivello, F., Etard, O., Delcroix, N., Mazoyer, B., and Joliot, M. (2002). Automated anatomical labeling of activations in SPM using a macroscopic anatomical parcellation of the MNI MRI single-subject brain. *NeuroImage*, 15(1):273–289. DOI: 10.1006/nimg.2001.0978.